# SYSTEM AND METHOD FOR AUTOMATIC RECOVERY FROM FAULT CONDITIONS IN NETWORKED COMPUTER SERVICES

### **CROSS-REFERENCE TO RELATED APPLICATION**

[0001] Not applicable.

# STATEMENT REGARDING FEDERALLY SPONSORED RESEARCH OR DEVELOPMENT

[0002] Not applicable.

### FIELD OF THE INVENTION

[0003] The invention relates to the field of network communications, and more particularly to a platform to monitor the overall performance and other operation of a computer network service, and to automatically deactivate or reroute defaulted services or connections to preserve service responsiveness.

#### **BACKGROUND OF THE INVENTION**

[0004] Deployment of networked computer services, such as Internet search engines, email, messaging and other communications platforms and others continue to expand and proliferate. Those and other services require Internet service providers (ISPs) and others to deploy increasingly capable back-end infrastructure to support the range and responsiveness of service expected by consumers, businesses and others. Installation of those resources, such as server farms, high-volume databases and others, in turn leads to demands for increased platform connectivity and results in greater dependency on all service components to cooperate effectively to deliver the search or other services.

[0005]

However, in an extensive installation such as a server farm arranged to support and Internet search engine or other application, the interdependency of numerous machines, connections and software may lead to faults or performance degradation in user-side performance when any one or more of the component resources crashes or becomes otherwise inoperable. For example a collection of servers which access travel, hotel and other remote data sources may hang or crash when executing a search on "Hawaii" or other terms when connections to one or more remote databases break or degrade. A user viewing a search page may therefore be presented with a blank screen, 404 error or other interruption of failure notification. This may occur even when other components, connections or data sources are still functioning and could perhaps return data to be presented to the user.

[0006]

In network service installations, to address that type of service interruption some operators may choose to install network monitoring packages which generate alerts to systems administrators, to advise them for example the processor utilization has become dangerously high on one group of servers, or that a backbone connection to a data source has broken down. This may permit the system administrator or other to step in and manually adjust communications links, activate redundant servers or take other actions. However, such arrangements still require the intervention and judgment of a human operator to sense and balance network performance in the presence of faults and other conditions. This among other things may lead to errors in judgment or a response time which is not acceptable or optimal during urgent network outages or conditions. Moreover human operators may only have the ability to monitor and act on a fairly limited number of connections or other resources for

emergency override purposes. Other problems in the management of networked computer services and the reliable operation of those services exist.

#### **SUMMARY OF THE INVENTION**

[0007]

The invention overcoming these and other problems in the art relates in one regard to a system and method for automatic recovery from fault conditions in networked computer services, in which one or more network monitors collect data on the performance of network servers, storage, connections and other resources and report network status data back to a control engine. In embodiments, the control engine may be configured with a rules-based logic engine which monitors the network performance data to detect and isolate emerging faults or other anomalous conditions, for example by detecting processor overload in a server or group of servers or a connection fault to a data source. Upon detecting that condition the control engine may automatically generate control commands to deactivate, disconnect or otherwise respond to the faulted server, connection or other resource. The end user or users accessing the affected resource, such as a search page, may consequently continue to view and receive results of the service, but with the failed component removed from the output or user interface. The control engine may likewise restore the faulted server, connection or other resources when the anomaly is removed, and may readjust other services to compensate for the affected service during the fault. According to embodiments of the invention in one regard, the resulting ability to monitor overall service health and "throw a switch" on one or more services when faults occur may occur may take place without necessary human intervention, although in implementations human operator overrides may be permitted.

#### **BRIEF DESCRIPTION OF THE DRAWINGS**

- [0008] Fig. 1 illustrates an network architecture in which a platform for automated recovery from a fault in networked computer services may operate, according to embodiments of the invention.
- [0009] Fig. 2 illustrates monitoring of network status data, according to embodiments of the invention.
- [0010] Fig. 3 illustrates dynamic feature or service adjustment, according to embodiments of the invention.
- [0011] Fig. 4 illustrates a flowchart of overall fault recovery processing, according to embodiments of the invention.

#### **DETAILED DESCRIPTION OF EMBODIMENTS**

[0012] Fig. 1 illustrates an architecture in which a system and method for automatic recovery from fault conditions in networked computer services may operate, according to an embodiment of the invention. As illustrated in that figure a consumer, corporate or other user may operate a user interface 102 to access a networked computer service 106 via a network 104, which may for example be, include or interface to the Internet, a local area network (LAN) or other network or connection. The service 106 to which the user connects via network 104 may be or include, for example, a search engine such as an Internet search engine which may search Web indices or other data stores for matching words, pictures or other content. Service 106 may likewise include communications or messaging services such as Microsoft

Network Messenger<sup>TM</sup>, an email client or service, a picture or other media exchange facility, or other applications or services.

In the environment as shown, in operation the user or users may be presented with a Web page or other service or content via user interface 102, which may be or include a browser or other client. The service or content presented via the user interface 102 may as shown include a set of features 106, such as linkable search results, email or other message links, advertising displays or applets, or other media, content or services. For example a user performing a search via a search engine may be presented with a set of configuration settings to control their search activity, for instance to restrict results to Web pages in English or other languages, or to permit the display of related or paid advertising sites.

[0014] Each feature in the set of features 106 may in turn interface to, communicate with or access other local or remote resources to support that feature. For example in a search engine application, a travel or vacation feature may connect to remote servers which serve selectable airfares, hotels or other related data or services. Other resources, such as databases or communications or messaging platforms, may likewise support and connect to the set of features 106. In embodiments, various of the set of features 106 may operate simultaneously, independently or in dependence on each other, the total interaction of which may contribute to the responsiveness and quality of the user's perceived experience when accessing service 106.

[0015] To manage and control potential faults, failures, degradations or interruptions in that content or other delivery, according to the invention in one regard a control

engine 108 may communicate with service 106 to monitor the health, status and performance of that service. Control engine 108 may be, include or be supported by server or other resources to maintain and dynamically adapt service 106 and supporting resources in the face of connection breakdowns, database outages and other faults and conditions. In configurations as shown, the control engine 108 may communicate with a control database 112, such as a relational database such as a structured query language (SQL) database or others, which in turn receives data from a network monitor 114.

[0016] Network monitor 114 may be or include a set of data capture ("sniffing") tools, connection monitoring or other functionalities which access and detect network status data 110 from any one or more components or services supporting the service 106, including network 104, connections to and from that network, processor and memory utilization in supporting servers or other machines, page latency in Web or other pages viewed by the user in user interface 102, or other parameters of network 104, service 106 or related resources. Network monitor 114 may for example be configured to monitor or "listen" to selected TCP/IP (transfer control protocol/Internet protocol) ports or other channels or interfaces to collect that network status data 110. That network status data 110 may be transmitted, streamed, sampled or stored to control database 112 for service management purposes.

[0017] As illustrated in Fig. 2, in embodiments the service 106 so managed may for example be or include a search engine, which service may be supported by a set of server, storage and other resources including a front end 118 for receiving search requests, a middle tier 120 for conditioning those requests, for example using SQL or

other query formats, and a backend 122 such as a database server or other data or analytic store. In embodiments involving a search engine application as part of service 106 as illustratively shown, the network monitor 114 may monitor the percentage of CPU utilization in front end 118, middle tier or back end 122 as part of the network status data 110, as shown. In embodiments a certain level of processor loading or other threshold or condition may be configured to trigger the detection of a system fault or other anomalous condition in control engine 108, such as 80% or other processor loading over a certain or sustained amount of time.

[0018]

When a fault, performance degradation or other condition is detected in this fashion, the control engine 108 may communicate control commands to the service 106 to dynamically adjust or alter the set of features 116 in response to the detected condition. For example if processor utilization has remained over threshold for a period of time in a server supporting a Web site which may normally be returned in search results, in embodiments the link to that site may be removed or deactivated to prevent users from attempting to link to that site while the exception or condition is taking place.

[0019]

As illustrated in Fig. 3, the control engine may for example dynamically adjust the set of features 116 of service 106 to ghost out or deactivate the normally available link for "www.siteD.com" due to processor loading at that site, connection difficulties with that site, or other conditions. By selectively deactivating that or other features in the set of features 116, the user may still be presented with useful information in user interface 102, so that the service 106 may continue while partial faults or degradations take place. Moreover since control engine 108 may detect and respond to one or more

network conditions on an automated basis, control actions may be taken on a realtime or near-realtime basis, which may be advantageous in cases of particularly severe or sudden faults. Moreover because the control engine 108 consolidates the monitoring of all points of support for service 106 in an integrated fashion, more than one of the set of features 116 or other service components may be controlled or deactivated at the same time in a coordinated fashion. That control oversight may also be configured to be conditional, for example to at first deactivate or suspend two or three of the set of features 116 when multiple or relatively widespread failures are being reported, to attempt to isolate the root or greatest contributing cause of the service fault or failure.

It may be also noted that in embodiments, the control commands and service adjustments generated by control engine 108 may be monitored or over-ridden by human administrators, when that configuration option is desired. Control engine 108 may in embodiments be configured to restore or release the deactivations or other adjustments to service 106, when the network status data 110 reflects that the fault, degradation or other condition has been alleviated. Decisioning on which resources to deactivate, reduce, adjust or otherwise manage in light of service or network conditions may in embodiments be supported or driven by rules-based logic in control engine 108, for example to adjust thresholds or other triggers depending on time of day, bandwidth utilization, user traffic and other variables. Other control options are possible.

[0021] Fig. 4 illustrates a flowchart of overall network monitoring and fault response processing, according to embodiments of the invention. In step 402, processing may

begin. In step 404, the status of network system parameters may be monitored, for instance via network monitor 114 or other access points. The captured network status data 110 may include for instance processor utilization rates, memory or storage usage, effective bandwidth, suspended or broken connections, page latency or other network or other conditions or parameters. In step 406, the network status data 110 reported by network monitor 114 or other tools or resources may be examined to determine whether network fault or other triggering conditions may exist, for example to determine whether a fault or other threshold has been met, for instance when processor utilization has reached 80% for an hour or other period, memory usage has instantaneously reached 90% or other levels, a backbone or other connection has been lost or degraded, or other triggering conditions or events are reported. The fault detection may in embodiments be performed by executing rules-based logic in control engine 108, for example by accessing control database 112 to load a threshold or group of related or conditional thresholds.

In step 408 when triggering conditions are detected, the control engine 114 may communicate a control command such as an interrupt to a faulted service or connection to deactivate, remove, suspend or otherwise manage the impact of that failure on the set of features 116 or other aspects of end-user experience or other performance variables. In step 410, the user interface 102 or content presented via that interface may for example be updated. For example in a search results page, an unavailable link may be ghosted out while still presenting other search results or other fields.

In step 412, the control engine 108 may further monitor and update the network status data and associated parameters, to determine any new or revised conditions arising from or related to the pending fault, or to detect unrelated anomalies. In step 414, the control engine 108 may determine whether release conditions have been met, so that the remedial actions taken in response to the current fault may be released. For example if a backbone connection has been verified to be restored, a total number of permitted users on an Internet site may be restored to two hundred from twenty, or other adjustments may be made or features in the set of features 116 reactivated. In step 416, when release or recovery conditions are met, the client interface 102 and related content may be dynamically restored or adjusted to reincorporate the restored service or feature in the set of features 116, or to make other links or resources available. In step 418, processing may repeat, return to a

[0024] The foregoing description of the invention is illustrative, and modifications in configuration and implementation will occur to persons skilled in the art. For instance, while the invention has generally been described in terms of a platform in which a singe control engine monitors all service components, in embodiments the logic and other functions amalgamated in the control engine could be deployed in multiple or distributed control servers or other engines. The control database likewise may in embodiments be distributed across multiple data stores.

prior processing point or terminate.

[0025] Similarly, while the invention has in embodiments been described as involving or supporting consumer-type networked services such as Internet search engines, other public or private, subscriber-based or non-subscriber based services or

applications may be monitored and managed according to embodiments of the invention. Other hardware, software or other resources described as singular may in embodiments be distributed, and similarly in embodiments resources described as distributed may be combined. The scope of the invention is accordingly intended to be limited only by the following claims.